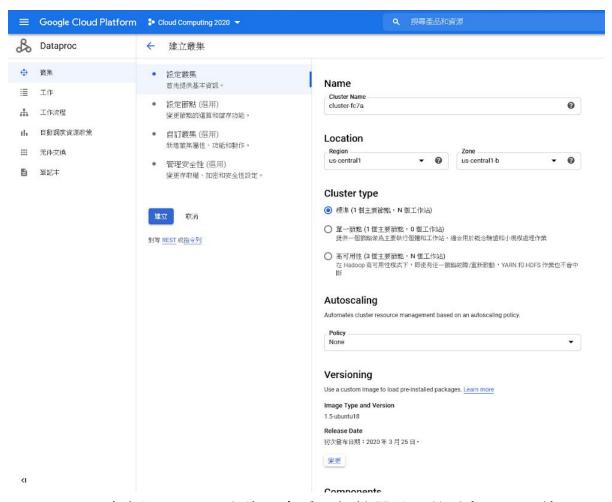
# **HW3 Hadoop Mapreduce**

r08942088 李政旻

- 1. (a) (90%) Write a map reduce application using Python to do the log analysis of apache2 web server access log
  - 1.1. 首先使用 google cloud platform 的 dataproc 來建立 cluster, 啟用 dataproc api 後,選擇標準的一個主要和預設兩個工作站的組合,由於配額的限制,所以兩個工作站都只使用一個CPU,而主要節點有4個CPU,並選擇ubuntu做作業系統。



1.2. 建立好 cluster 之後,查看三個機器分別的外部IP,再使用 SSH進入機器確認環境是否正常。



```
A520M-4750G@cluster-d98c-m:~$ hadoop envvars

JAVA_HOME='/usr/lib/jvm/adoptopenjdk-8-hotspot-amd64'

HADOOP_COMMON_HOME='/usr/lib/hadoop'

HADOOP_COMMON_DIR='./'

HADOOP_COMMON_LIB_JARS_DIR='lib'

HADOOP_COMMON_LIB_NATIVE_DIR='lib/native'

HADOOP_CONF_DIR='/etc/hadoop/conf'

HADOOP_TOOLS_PATH='/usr/lib/hadoop/share/hadoop/tools/lib/*'
```

- 1.3. 使用SFTP將需要的檔案上傳 (access.log, mapper.py, reducer.py, test.sh)
  - 1.3.1. access.log: 即
    <a href="http://hpc.ee.ntu.edu.tw/html/IntelligentClouds/webAccessLog/accesslog">http://hpc.ee.ntu.edu.tw/html/IntelligentClouds/webAccessLog/accesslog</a>
  - 1.3.2. mapper.py:

```
#!/usr/bin/env python
      ""mapper.py"
     import sys
     import datetime
     # input comes from STDIN (standard input)
    pfor line in sys.stdin:
          # remove leading and trailing whitespace
         line = line.strip()
12
         # extract record time
         words = line.split("[")
14
         if len(words) > 1:
15
             record time = words[1].split()
         else:
16
             continue
18
19
         # change time format
         formated record time = datetime.datetime.strptime(record time[0], "%d/%b/%Y:%H:%M:%S")
20
21
         print (formated record time.strftime('%Y-%m-%d T %H:00:00.000')+"\t1")
```

## 1.3.3. reducer.py:

```
#!/usr/bin/env python
     """reducer.py"""
 3
 4
     from operator import itemgetter
 5
     import sys
 6
 7
     current word = None
 8
     current count = 0
 9
     word = None
10
11
     # input comes from STDIN
12
    pfor line in sys.stdin:
         # remove leading and trailing whitespace
13
14
         line = line.strip()
15
16
         # parse the input we got from mapper.py
17
        word, count = line.split('\t', 1)
18
19
         # convert count (currently a string) to int
20 白
         try:
21
             count = int(count)
22 白
         except ValueError:
23
             # count was not a number, so silently
24
             # ignore/discard this line
25
             continue
26
27
         # this IF-switch only works because Hadoop sorts map output
28
         # by key (here: word) before it is passed to the reducer
29
         if current word == word:
             current count += count
31
         else:
32
             if current word:
33
                 # write result to STDOUT
34
                 print ('%s\t%s' % (current_word, current_count))
35
             current count = count
36
             current word = word
37
    # do not forget to output the last word if needed!
39 ⊟if current word == word:
40
         print ('%s\t%s' % (current word, current count))
41
```

#### 1.3.4. test.sh:

#### 1.4. 執行 "bash test.sh"

```
20/11/03 18:18:27 INFO mapreduce.Job: map 80% reduce 0% 20/11/03 18:18:33 INFO mapreduce.Job: map 87% reduce 0%
20/11/03 18:18:36 INFO mapreduce. Job: map 100% reduce 0%
20/11/03 18:18:44 INFO mapreduce.Job: map 100% reduce 100% 20/11/03 18:18:47 INFO mapreduce.Job: Job job_1604426318533_0001 completed successfully
20/11/03 18:18:48 INFO mapreduce. Job: Counters: 49
        File System Counters
                 FILE: Number of bytes read=46392
                 FILE: Number of bytes written=3824095
                 FILE: Number of read operations=0
FILE: Number of large read operations=0
                 FILE: Number of write operations=0
                 HDFS: Number of bytes read=233143
                 HDFS: Number of bytes written=2825
                 HDFS: Number of read operations=55
                 HDFS: Number of large read operations=0
                 HDFS: Number of write operations=6
         Job Counters
                 Launched map tasks=15
                 Launched reduce tasks=2
                 Data-local map tasks=15
                 Total time spent by all maps in occupied slots (ms)=784232
                 Total time spent by all reduces in occupied slots (ms)=87640
                 Total time spent by all map tasks (ms)=196058
                 Total time spent by all reduce tasks (ms)=10955
                 Total vcore-milliseconds taken by all map tasks=196058
                 Total vcore-milliseconds taken by all reduce tasks=10955
                 Total megabyte-milliseconds taken by all map tasks=200763392
                 Total megabyte-milliseconds taken by all reduce tasks=22435840
        Map-Reduce Framework
                 Map input records=1546
                 Map output records=1546
                 Map output bytes=43288
                 Map output materialized bytes=46560
                 Input split bytes=1350
                 Combine input records=0
                 Combine output records=0
                 Reduce input groups=99
                 Reduce shuffle bytes=46560
                 Reduce input records=1546
                 Reduce output records=99
                 Spilled Records=3092
                 Shuffled Maps =30
                 Failed Shuffles=0
                 Merged Map outputs=30
                 GC time elapsed (ms)=4007
                 CPU time spent (ms)=12700
                 Physical memory (bytes) snapshot=6725951488
Virtual memory (bytes) snapshot=45217533952
                 Total committed heap usage (bytes) = 5027065856
        Shuffle Errors
                 BAD ID=0
                 CONNECTION=0
                 IO ERROR=0
                 WRONG LENGTH=0
                 WRONG MAP=0
                 WRONG REDUCE=0
        File Input Format Counters
                 Bytes Read=231793
         File Output Format Counters
                 Bytes Written=2825
20/11/03 18:18:48 INFO streaming.StreamJob: Output directory: /output
A520M-4750G@cluster-d98c-m:~/hw3$
```

#### 1.5. 執行 "cat result.dat" 查看結果

```
A520M-4750G@cluster-d98c-m:~/hw3S cat result.dat
2004-03-07 T 16:00:00.000
                                27
2004-03-07 T 17:00:00.000
                                25
2004-03-07 T 18:00:00.000
                                24
2004-03-07 T 19:00:00.000
                                26
2004-03-07 T 20:00:00.000
                                20
2004-03-07 T 21:00:00.000
                                23
2004-03-07 T 22:00:00.000
                                29
2004-03-07 T 23:00:00.000
                                22
2004-03-08 T 00:00:00.000
                                21
2004-03-08 T 01:00:00.000
                                21
2004-03-08 T 02:00:00.000
                                27
2004-03-08 T 03:00:00.000
                                22
2004-03-08 T 04:00:00.000
                                26
2004-03-08 T 05:00:00.000
                                37
2004-03-08 T 06:00:00.000
                                17
2004-03-08 T 07:00:00.000
                                31
2004-03-08 T 08:00:00.000
                                44
2004-03-08 T 09:00:00.000
                                63
2004-03-08 T 10:00:00.000
                                39
2004-03-08 T 11:00:00.000
                                34
2004-03-08 T 12:00:00.000
                                45
2004-03-08 T 13:00:00.000
                                37
2004-03-08 T 14:00:00.000
                                23
2004-03-08 T 15:00:00.000
                                9
2004-03-08 T 16:00:00.000
                                2
2004-03-08 T 17:00:00.000
                                2
2004-03-08 T 18:00:00.000
                                9
2004-03-08 T 19:00:00.000
                                6
2004-03-08 T 20:00:00.000
                                23
2004-03-08 T 22:00:00.000
                                20
2004-03-08 T 23:00:00.000
2004-03-09 T 01:00:00.000
                                12
```

#### 1.6. reference:

- 1.6.1. <a href="https://ierrynest.io/dataproc-hadoop/">https://ierrynest.io/dataproc-hadoop/</a>
- 1.6.2. <a href="https://ierrvnest.io/install-hadoop-on-ubuntu/">https://ierrvnest.io/install-hadoop-on-ubuntu/</a>

- 1.6.3. <a href="https://jerrynest.io/hadoop-map-reduce-apache2-web-s">https://jerrynest.io/hadoop-map-reduce-apache2-web-s</a> erver-access-log/
- 2. (b) (10%) Write the Java version of part (a)

### 2.1. hourCount.java:

```
Epublic class hourCount {
19
       public static class TokenizerMapper
            extends Mapper<Object, Text, Text, IntWritable>{
         private final static IntWritable one = new IntWritable(1);
23
         private Text word = new Text();
24
25
         public void map (Object key, Text value, Context context
26
                         ) throws IOException, InterruptedException (
27
            /*StringTokenizer itr = new StringTokenizer(value.toString());
28
           while (itr.hasMoreTokens()) {
29
             word.set(itr.nextToken());
             context.write(word, one);
31
           String[] temp = value.toString().split("\\[");
           if (temp.length > 1) {
34
             temp = temp[1].split("\\s+");
             SimpleDateFormat sdf_input = new SimpleDateFormat("dd/MMM/yyyy:HH:mm:ss");
36
37
             SimpleDateFormat sdf_output = new SimpleDateFormat("yyyy-MM-dd HH:00:00.000");
39
             try {
                 word.set(sdf_output.format(sdf_input.parse(temp[0])));
40
                 context.write (word, one);
41
42
              } catch (ParseException e) {
43
44
           }
45
46
47
       }
48
```

```
49
       public static class IntSumReducer
50 白
            extends Reducer<Text, IntWritable, Text, IntWritable> {
         private IntWritable result = new IntWritable();
51
52
53
        public void reduce (Text key, Iterable < IntWritable > values,
54
                             Context context
55 白
                             ) throws IOException, InterruptedException {
56
          int sum = 0:
57 白
          for (IntWritable val : values) {
58
             sum += val.get();
59
60
           result.set(sum);
61
           context.write(key, result);
62
63
       }
64
65 public static void main (String[] args) throws Exception {
         Configuration conf = new Configuration();
66
67
         Job job = Job.getInstance(conf, "word count");
68
         job.setJar("hc.jar");
69
         job.setJarByClass(hourCount.class);
70
         job.setMapperClass(TokenizerMapper.class);
71
         job.setCombinerClass(IntSumReducer.class);
72
         job.setReducerClass(IntSumReducer.class);
73
         job.setOutputKeyClass(Text.class);
74
         job.setOutputValueClass(IntWritable.class);
75
         FileInputFormat.addInputPath(job, new Path(args[0]));
76
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
77
         System.exit(job.waitForCompletion(true) ? 0 : 1);
78
79
```

#### 2.2. 送上 cluster 後. 依序執行:

- 2.2.1. hdfs dfs -mkdir /input
- 2.2.2. hdfs dfs -put access.log /input
- 2.2.3. hdfs dfs -rm -r -f /output
- 2.2.4. javac hourCount.java -cp \$(hadoop classpath)
- 2.2.5. jar cf hc.jar hourCount\*.class
- 2.2.6. hadoop jar hc.jar hourCount /input /output
- 2.2.7. hdfs dfs -cat /output/\* | sort -k1,1

```
A520M-4750G@cluster-d98c-m:~/hw3/java_hour_count$ hdfs dfs -cat /output/* | sort -k1,1
2004-03-07 16:00:00.000 27
2004-03-07 17:00:00.000 25
2004-03-07 18:00:00.000 24
2004-03-07 19:00:00.000 26
2004-03-07 20:00:00.000 20
2004-03-07 21:00:00.000 23
2004-03-07 22:00:00.000 29
2004-03-07 23:00:00.000 22
2004-03-08 00:00:00.000 21
2004-03-08 01:00:00.000 21
2004-03-08 02:00:00.000 27
2004-03-08 03:00:00.000 22
2004-03-08 04:00:00.000 26
2004-03-08 05:00:00.000 37
2004-03-08 06:00:00.000 17
2004-03-08 07:00:00.000 31
2004-03-08 08:00:00.000 44
2004-03-08 09:00:00.000 63
2004-03-08 10:00:00.000 39
2004-03-08 11:00:00.000 34
2004-03-08 12:00:00.000 45
2004-03-08 13:00:00.000 37
2004-03-08 14:00:00.000 23
2004-03-08 15:00:00.000 9
2004-03-08 16:00:00.000 2
2004-03-08 17:00:00.000
2004-03-08 18:00:00.000 9
2004-03-08 19:00:00.000 6
2004-03-08 20:00:00.000 23
2004-03-08 22:00:00.000 20
2004-03-08 23:00:00.000 1
2004-03-09 01:00:00.000 12
2004-03-09 02:00:00.000 15
2004-03-09 03:00:00.000 1
```

2.3. reference: <a href="https://jerrynest.io/hadoop-wordcount-example/">https://jerrynest.io/hadoop-wordcount-example/</a>

## 3. github url:

https://github.com/a2134666/CloudComputingHW3